

VC dimension  
and  
Model Selection

# Overview

□ PAC model: review

□ VC dimension:

➤ Definition

➤ Examples

➤ Sample:

- Lower bound
- Upper bound !!!

□ Model Selection

# PAC model: Setting

- A distribution:  $D$  (unknown)
- Target function:  $c_t$  from  $C$ 
  - $c_t : X \rightarrow \{0,1\}$
- Hypothesis:  $h$  from  $H$ 
  - $h : X \rightarrow \{0,1\}$
- Error probability:
  - $\text{error}(h) = \text{Prob}_D[h(x) \neq c_t(x)]$
- Oracle:  $EX(c_t, D)$

# PAC model: Definition

- C and H are concept classes over X.
- C is PAC learnable by H if
- There Exist an Algorithm A such that:
  - For any distribution D over X and  $c_t$  in C
  - for every input  $\epsilon$  and  $\delta$ :
  - outputs a hypothesis h in H,
    - while having access to  $EX(c_t, D)$
  - with probability  $1-\delta$  we have  $\text{error}(h) < \epsilon$
- Complexities: sample, running time

# PAC model – last week

□ For a finite hypothesis class  $H$ , sample size  $m$ :

➤ Realizable case:

$$m > (1/\varepsilon) \ln (|H|/\delta)$$

➤ Non-realizable

$$m > (2/\varepsilon^2) \ln (2|H|/\delta)$$

□ Impossibility results:

$$➤ m > (1/2) \log |H|$$

$$➤ m > (1/4\varepsilon) \ln (1/\delta)$$

# VC dimension: motivation

## □ Infinite hypothesis class

➤ Threshold

➤ Rectangles

➤ **TODAY: general**

- VC dimension

- Applies both to realizable and non-realizable.

# VC dimension: definition

- Notation:  $C$  – concept class;  $S$  - sample
- Projection:  $\Pi_C(S) = \{c \cap S : c \in C\}$
- Shattering:  $C$  shatters  $S$  if  $|\Pi_C(S)| = 2^{|S|}$
- VC dimension: size of largest  $S$  shattered
  - $\max\{d: \exists S, |S| = d, \Pi_C(S) = 2^{|S|}\}$
  - If “no max” then infinity
    - For every  $d$  there is a shattered set of size  $d$

# VC dimension: Threshold

□  $c_\theta(x) = I(x \geq \theta)$



□  $VC \geq 1$

➤  $S = \{0.5\}$ :  $c_{0.3}(0.5) = 1$  and  $c_{0.6}(0.5) = 0$

□  $VC < 2$

➤  $S = \{z_1, z_2\}$

➤ Assume  $z_1 < z_2$

➤  $c(z_1) = 1, c(z_2) = 0$



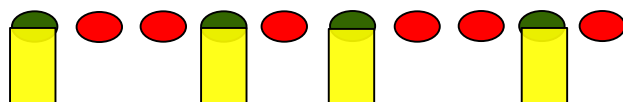
# VC dimension: union of intervals

□ Intervals on  $[0,1]$



➤ Finite but unbounded

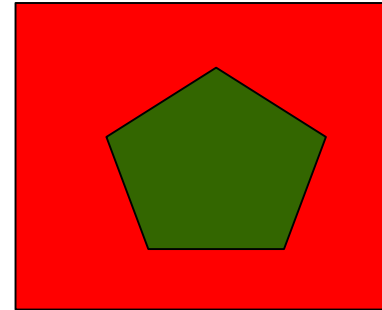
□ For any  $d$  points:



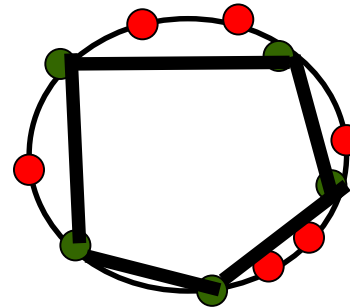
□ VC-dim = infinity

# VC dimension: convex polygon

□ Convex polygon



□ For any  $d$  points



□ VC dimension = infinity

# VC dimension: hyperplane

$$\square c_{w,\theta}(x) = \text{sign}(\sum_i w_i x_i + \theta)$$

$$\square \text{VC dimension} \geq d+1$$

$$\triangleright S = \{\vec{0}, \vec{e}_1, \dots, \vec{e}_d\}$$

$\triangleright$  Given labeling  $L$  in  $\{-1, +1\}$ , define  $c_{w,\theta}(x)$ :

$$\bullet w_i = L(x_i)$$

$$\bullet \theta = \frac{L(0)}{2}$$

$$\triangleright c_{w,\theta}(0) = \text{sign}(\theta) \ \& \ c_{w,\theta}(e_i) = \text{sign}(w_i + \theta)$$

# VC dimension: hyperplane

□ VC dimension  $< d+2$

➤ For contradiction

- Assume there is a shattered  $S$  and  $|S|=d+2$

□ Radom Theorem:  $S \subseteq \mathbb{R}^d$ ,  $|S| \geq d + 2$

➤  $\exists S' \subset S : \text{conv}(S') \cap \text{conv}(S - S') \neq \emptyset$

□ Let  $S'$  be positive and  $S-S'$  be negative

➤ Let  $c_{w,\theta}$  be the separating hyperplane

- Let POS be the positive and NEG be negatives of  $c_{w,\theta}$

# VC dimension: hyperplane

- $\text{conv}(S') \subset \text{POS}$  &  $\text{conv}(S - S') \subset \text{NEG}$ 
  - closed under convex combinations
- Radom Theorem:
  - $\text{conv}(S') \cap \text{conv}(S - S') \neq \emptyset$
- However:  $\text{POS} \cap \text{NEG} = \emptyset$
- Contradiction!
  - There is no such set  $S$
- $\text{VC-dim} < d+2$

□ QED

# VC dim: Sample lower bound

□ Theorem:

$$VC\text{-dim}(C) = d + 1$$

$$\triangleright m \geq \frac{d}{16\epsilon}$$

□ Proof:

➤ Let  $\{z_0, z_1, \dots, z_d\}$

➤  $D(x)$

$$1 - 8\epsilon \quad x = z_0$$

$$= \frac{8\epsilon}{d} \quad x = z_i$$

$$0 \quad \text{otherwise}$$

□ Target function:

$$\triangleright c_t(z_0) = 1;$$

$$\triangleright c_t(z_i) = 0 \text{ or } 1 \text{ (prob } \frac{1}{2})$$

□  $RARE = \{z_1, \dots, z_d\}$

$$\triangleright \text{Assume } |S \cap RARE| \leq \frac{d}{2}$$

$$\triangleright |UNSEEN| \geq \frac{d}{2}$$

$$\triangleright \Pr[\text{error}] \geq \frac{1}{2} \frac{8\epsilon}{d} |UNSEEN| \geq 2\epsilon$$

# VC dim: Sample lower bound

$$\square E[|S \cap RARE|] = 8\epsilon m \leq \frac{d}{2}$$

$$\square \Pr \left[ |S \cap RARE| \leq \frac{d}{2} \right] \geq \frac{1}{2}$$

$\square$  With probability at least  $\frac{1}{2}$

$\blacktriangleright$  Error at least  $2\epsilon$

$\square$  QED

# VC dim: sample upper bound

❑ Incorrect proof

❑ For sample  $S$ :

➤  $C_{|S} = \Pi_C(S)$  is finite

❑ Use finite class bound:

➤  $m \geq \frac{1}{\epsilon} \log \frac{|\Pi_C(S)|}{\delta}$

❑ Problem:

➤  $S$  defines  $C_{|S} = \Pi_C(S)$

❑ Solution

➤ Take  $2m$  points

- $S = S_1 \cup S_2$
- The randomization in the split to  $S_1$  and  $S_2$

➤ **Benefit:**

- We have  $\Pi_C(S)$



# VC dim: sample upper bound

## □ Bad concepts

- $\text{Bad} = \{h \mid \text{error}(h) > \epsilon\}$

## □ Hitting set $S$ :

- For every  $h$  in  $\text{Bad}$
- Exists  $x$  in  $S$
- $c_t(x) \neq h(x)$

## □ Goal

- Compute prob. of  $S$  being a hitting set

## □ Event A:

- $S_1$  not hitting set
  - Exists  $h$  in  $\text{Bad}$  which is consistent
- $\Pr[A] < ???$

## □ Event B:

- Exists  $h$  in  $\text{Bad}$ 
  - $h$  consistent with  $S_1$
  - $h$  has  $\epsilon m$  errors on  $S_2$

# VC dim: sample upper bound

$$\square \Pr[B] = \Pr[B|A] \Pr[A]$$

➤ Since B implies A

$$\square \Pr[B|A]$$

➤ Fix such an  $h$

➤ Expected errors  $\geq \epsilon m$

➤ Probability at least  $\frac{1}{2}$

$\square$  Result:

➤  $2 \Pr[B] \geq \Pr[A]$

$$\square F = \Pi_C(S_1 \cup S_2)$$

$\square$  Fix  $h$  in  $F$ :

➤  $h$  consistent with  $S_1$

➤  $h$  has errors  $\geq \epsilon m$  on  $S_2$

•  $l$  number of errors

$\square$  Compute the prob.  
over partitions

➤  $S_1$  and  $S_2$

# VC dim: sample upper bound

□ Number of total partitions:

$$\triangleright \binom{2m}{l}$$

□ Number of partitions which make  $h$  consistent on  $S_1$

$$\triangleright \binom{m}{l}$$

□ Prob bound

$$\triangleright \frac{\binom{m}{l}}{\binom{2m}{l}} = \prod_{i=0}^{l-1} \frac{m-i}{2m-i} \leq \frac{1}{2^l}$$

□ Bounding probabilities:

‣ Union bound over  $h$  in  $F$

$$\triangleright \Pr[B] \leq |F|2^{-\varepsilon m}$$

$$\triangleright \Pr[A] \leq 2\Pr[B] \leq 2|F|2^{-\varepsilon m}$$

# VC dim: sample upper bound

## □ High confidence

➤  $\delta \geq 2|F|2^{-\epsilon m}$

➤  $m \geq \frac{1}{\epsilon} \log \frac{2|F|}{\delta}$

## □ Need to bound $|F|$

➤  $F = \Pi_C(S_1 \cup S_2)$

## □ Sauer-Shelah Lemma:

➤  $VC\text{-dim}(C)=d$

➤  $|S|=2m$

➤  $\Pi_C(S) \leq \sum_{i=0}^d \binom{2m}{i}$

## □ Bound

➤  $2^m$  for  $m \leq d$

➤  $2(2m)^d$  for  $m > d$

# VC dim: Sampling Theorem

## □ Sample bound

- $m \geq \frac{1}{\epsilon} \log \frac{4(2m)^d}{\delta}$
- $m \geq \frac{2}{\epsilon} + \frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{d}{\epsilon} \log 2m$
- $m = O\left(\frac{1}{\epsilon} \log \frac{1}{\delta} + \frac{d}{\epsilon} \log \frac{d}{\epsilon}\right)$

## □ Realizable case

## □ Non-realizable

- $m = O\left(\frac{1}{\epsilon^2} \log \frac{1}{\delta} + \frac{d}{\epsilon^2} \log \frac{d}{\epsilon}\right)$

## □ Proof methodology

# Rademacher Complexity

## □ Motivation:

- Tighter bounds; Dist. Dependent

## □ Notation:

- $f \in \{-1, +1\}; f \in F$

- $\Pr[\sigma_i = +1] = \Pr[\sigma_i = -1] = \frac{1}{2}$

# Rademacher Complexity

## □ Definition (Rademacher Complexity):

➤  $S$  sample of size  $m$

$$\text{➤ } R_S(F) = E_{\sigma} \left[ \max_{f \in F} \sum_{i=1}^m \sigma_i f(x_i) \right]$$

$$\text{➤ } R_D(F) = E_S [R_S(F)]$$

# Rademacher Complexity: expected overfitting

□ Theorem (expected overfitting):

$$\blacktriangleright E_S \left[ \max_{f \in F} \frac{1}{m} \sum_{i=1}^m f(x_i) - E_D[f(x)] \right] \leq 2R_D(F)$$

□ Proof:

➤ Two sample trick, add  $S'$

$$\blacktriangleright = E_S \left[ \max_{f \in F} \frac{1}{m} \sum_{i=1}^m f(x_i) - E_{S'} \left[ \frac{1}{m} \sum_{i=1}^m f(x'_i) \right] \right]$$

$$\blacktriangleright \leq E_{S, S'} \left[ \max_{f \in F} \frac{1}{m} \sum_{i=1}^m f(x_i) - f(x'_i) \right]$$



# Rademacher Complexity: expected overfitting

$$\square = E_{S, S'} \left[ \max_{f \in F} \frac{1}{m} \sum_{i=1}^m \sigma_i (f(x_i) - f(x'_i)) \right]$$

$$\square \leq E_S \left[ \max_{f \in F} \frac{1}{m} \sum_{i=1}^m \sigma_i f(x_i) \right] + E_{S'} \left[ \max_{f \in F} \frac{1}{m} \sum_{i=1}^m \sigma_i f(x'_i) \right]$$

$$\square = 2R_D(F)$$

□ QED

# Rademacher Theorem

□ With probability  $1-\delta$ , for every  $h \in H$ :

$$\begin{aligned} \blacktriangleright \epsilon(h) &\leq \hat{\epsilon}(h) + R_D(H) + \sqrt{\frac{\ln\left(\frac{2}{\delta}\right)}{2m}} \\ &\leq \hat{\epsilon}(h) + R_S(H) + 3\sqrt{\frac{\ln\left(\frac{2}{\delta}\right)}{2m}} \end{aligned}$$

# Model selection - Outline

- Motivation
- Overfitting
- Structural Risk Minimization
- Hypothesis Validation

# Motivation:

## □ Problems:

- We have too few examples
- We have a very rich hypothesis class
- How can we find the best hypothesis?

## □ Alternatively:

- Usually we choose the hypothesis class
  - How rich of a class we want?
- ## □ How should we go about doing it?

# Overfitting

- ❑ Concept class: Intervals on a line
- ❑ Can classify any training set
- ❑ Zero training error:
  - Is this the only goal?!



# Overfitting: Intervals



- ❑ Can always get zero training error!
- ❑ Are we interested in zero training error?!

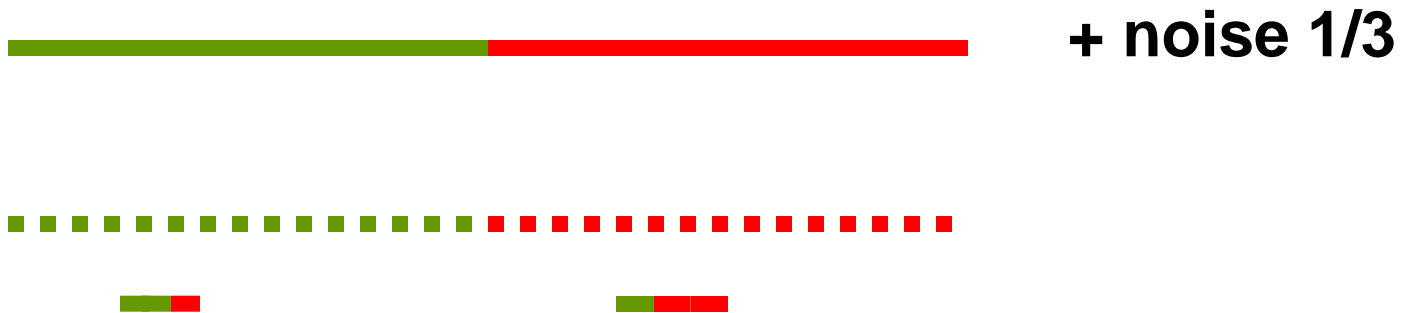
# Overfitting: Intervals



intervals	0	1	2	3	4
errors	7	3	2	1	0

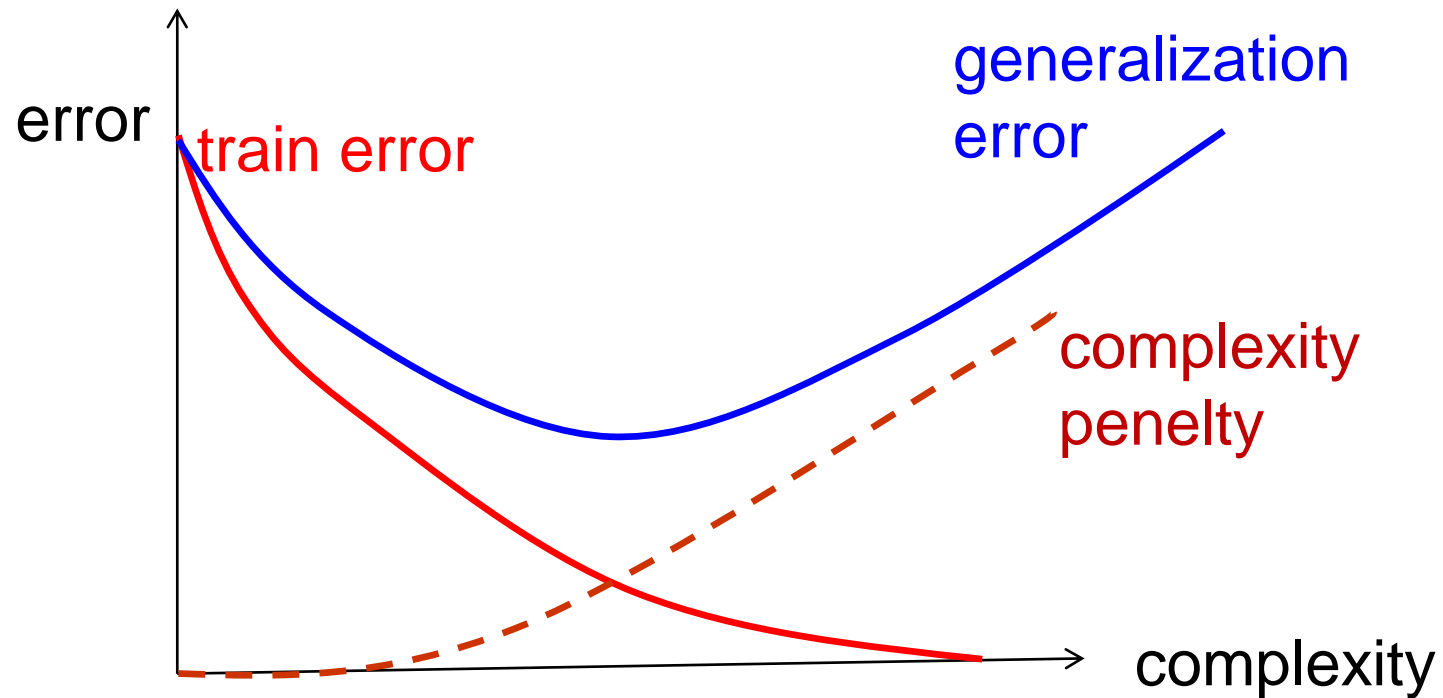
# Overfitting

- ❑ Simple concept plus noise
- ❑ A very complex concept
  - insufficient number of examples





# Model Selection



# Theoretical Model

□ Nested Hypothesis classes

➤  $H_1 \subseteq H_2 \subseteq H_3 \subseteq \dots \subseteq H_i \subseteq$

□ There is a target function  $c_t(x)$ ,

➤ non-realizable.

□ True errors:

➤  $\varepsilon(h) = Pr [ h \neq c_t ]$

➤  $\varepsilon_i = \mathbf{inf}_{h \in H_i} e(h)$

➤  $\varepsilon(h^*) = \mathbf{inf}_i \varepsilon_i$

•  $h^*$  is best hypothesis

□ Training error

➤  $\hat{\varepsilon}(h) = \frac{1}{m} \sum_{i=1}^m I[ h \neq c_t ]$

➤  $\hat{\varepsilon}_i = \mathbf{inf}_{h \in H_i} \hat{\varepsilon}(h)$

# Theoretical Model

## □ Complexity of $h$

➤  $d(h) = \mathbf{\min}_i \{h \in H_i\}$

## □ Add a penalty for $d(h)$

➤ minimize:  $\hat{e}(h) + \text{penalty}(h)$

## □ Penalty based.

➤ Chose the hypothesis which minimizes:

➤  $\hat{e}(h) + \text{penalty}(h)$

# Structural Risk Minimization

□ Parameters:  $\lambda_i$  and  $\delta_i$  such that:

$$\Pr[\exists h \in H_i: |\hat{\epsilon}(h) - \epsilon(h)| > \lambda_i] \leq \delta_i$$

➤  $\sum_i \delta_i = \delta$

•  $\delta_i = \delta/2^i$

□ Implies: with prob.  $1-\delta$

➤  $\Pr[\exists h \in H: |\hat{\epsilon}(h) - \epsilon(h)| > \lambda_{d(h)}] \leq \delta_{d(h)}$

# Structural Risk Minimization

□ Setting  $penalty(h) = \lambda_{d(h)}$

□ Finite  $|H_i|$

$$\triangleright \lambda_i = \sqrt{\frac{\log |H_i| / \delta}{m}}$$

□  $VC-dim(H_i) = i$

$$\triangleright \lambda_i = \sqrt{\frac{i \log i / \delta}{m}}$$

# SRM: Performance

## □ THEOREM

- $h^*$  : best hypothesis
- $g_{srm}$  : SRM choice
- With probability  $1-\delta$
- $\varepsilon(h^*) \leq \varepsilon(g_{srm}) \leq \varepsilon(h^*) + 2 \text{penalty}(h^*)$

□ Note: bound depends only on  $h^*$

# Proof

□ Bounding the error in  $H_i$

$$\begin{aligned} &\triangleright \Pr[|\hat{\epsilon}(g_{srm}) - \epsilon(g_{srm})| > \lambda_{srm}] \\ &\leq \Pr[\exists h \in H_{srm}: |\hat{\epsilon}(h) - \epsilon(h)| > \lambda_{srm}] \leq \delta_{srm} \end{aligned}$$

□ Bounding the error across  $H_i$

$$\begin{aligned} &\triangleright \hat{\epsilon}(g_{srm}) \geq \epsilon(g_{srm}) - \lambda_{srm} \\ &\triangleright \hat{\epsilon}(h^*) + \lambda_* \geq \hat{\epsilon}(g_{srm}) + \lambda_{srm} \\ &\triangleright \epsilon(h^*) + \lambda_* \geq \hat{\epsilon}(h^*) \end{aligned}$$

$$\square \epsilon(h^*) + 2\lambda_* \geq \epsilon(g_{srm}) \quad \text{QED}$$

# Hypothesis Validation

- ❑ Separate sample to training and selection.

- ❑ Using the training

  - Select from each  $H_i$  a candidate  $g_i$

- ❑ Using the selection sample

  - select between  $g_1, \dots, g_m$

- ❑ The split size

  - $(1-\gamma)m$  training set

  - $\gamma m$  selection set



# Hypothesis Validation: Algorithm

□ Using  $(1-\gamma)m$  examples:  $S_1$

➤  $\hat{\epsilon}_1(h)$  = error on  $S_1$

➤  $g_i = \arg \min_{h \in H_i} \hat{\epsilon}_1(h)$

□ Using  $\gamma m$  examples:  $S_2$

➤  $\hat{\epsilon}_2(h)$  = error on  $S_2$

➤  $g_{HV} = \arg \min_{g_i \in G} \hat{\epsilon}_2(g_i)$

□ Return  $g_{HV}$

# Hypo. Validation: Performance

## □ Errors

- $\varepsilon_{hv}(m)$  = error of HV
  - Using  $m$  examples
- $\varepsilon_A(m)$  = error of A
  - Any algorithm
  - Using  $m$  examples
  - Selecting  $g_i$  from  $H_i$ 
    - only restriction on A
  - For example: any penalty function

□ Theorem: with probability  $1-\delta$

$$\varepsilon_{hv}(m) \leq \varepsilon_A((1-\gamma)m) + 2\sqrt{\frac{\ln(2m/\delta)}{\gamma m}}$$

# Hypo. Validation: Analysis

$$\square \Pr[|\epsilon(g_i) - \hat{\epsilon}_2(g_i)| > \lambda] \leq 2e^{-\lambda^2 \gamma m}$$

$$\square \Pr[\exists i: |\epsilon(g_i) - \hat{\epsilon}_2(g_i)| > \lambda] \leq 2|G|e^{-\lambda^2 \gamma m} = \delta$$

$$\square \text{Since } |G| \leq m: \lambda = \sqrt{\frac{\ln 2m/\delta}{\gamma m}}$$

$$\triangleright \epsilon_2(g_i) + \lambda \geq \hat{\epsilon}_2(g_i)$$

$$\triangleright \hat{\epsilon}_2(g_i) \geq \hat{\epsilon}_2(g_{HV})$$

$$\triangleright \hat{\epsilon}_2(g_{HV}) \geq \epsilon(g_{HV}) - \lambda$$

$$\square \epsilon_2(g_i) + 2\lambda \geq \epsilon(g_{HV})$$

# Summary

- PAC model

- Generalization bounds

  - Empirical Risk Minimization

  - VC dimension

  - Rademacher complexity

- Model Selection

  - Structural Risk Minimization (SRM)

  - Hypothesis selection